



On Calibrating Diffusion Probabilistic Models

Zhijie Deng

zhijied@sjtu.edu.cn

Shanghai Jiao Tong University

Joint work with Tianyu Pang, Cheng Lu, Chao Du, Min Lin, and Shuicheng Yan

Diffusion Models in 2020 (Nonequilibrium Thermodynamics)

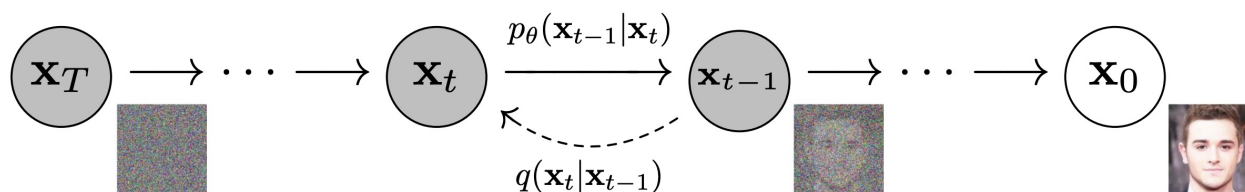


Figure 2: The directed graphical model considered in this work.

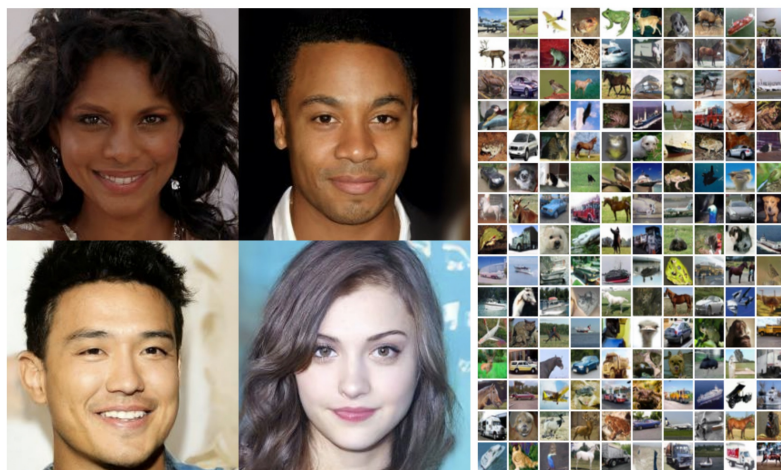


Figure 1: Generated samples on CeleBA-HQ 256 x 256 (left) and unconditional CIFAR10 (right)



Figure 3: LSUN Church samples. FID=7.89

Figure 4: LSUN Bedroom samples. FID=4.90

- [1] Sohl-Dickstein et al. Deep Unsupervised Learning using Nonequilibrium Thermodynamics. ICML 2015
- [2] Ho et al. Denoising Diffusion Probabilistic Models. NeurIPS 2020

Diffusion Models in 2020 (Annealed Langevin Dynamics)



Algorithm 1 Annealed Langevin dynamics.

Require: $\{\sigma_i\}_{i=1}^L, \epsilon, T$.

- 1: Initialize $\tilde{\mathbf{x}}_0$
 - 2: **for** $i \leftarrow 1$ to L **do**
 - 3: $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ $\triangleright \alpha_i$ is the step size.
 - 4: **for** $t \leftarrow 1$ to T **do**
 - 5: Draw $\mathbf{z}_t \sim \mathcal{N}(0, I)$
 - 6: $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
 - 7: **end for**
 - 8: $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$
 - 9: **end for**
- return** $\tilde{\mathbf{x}}_T$
-

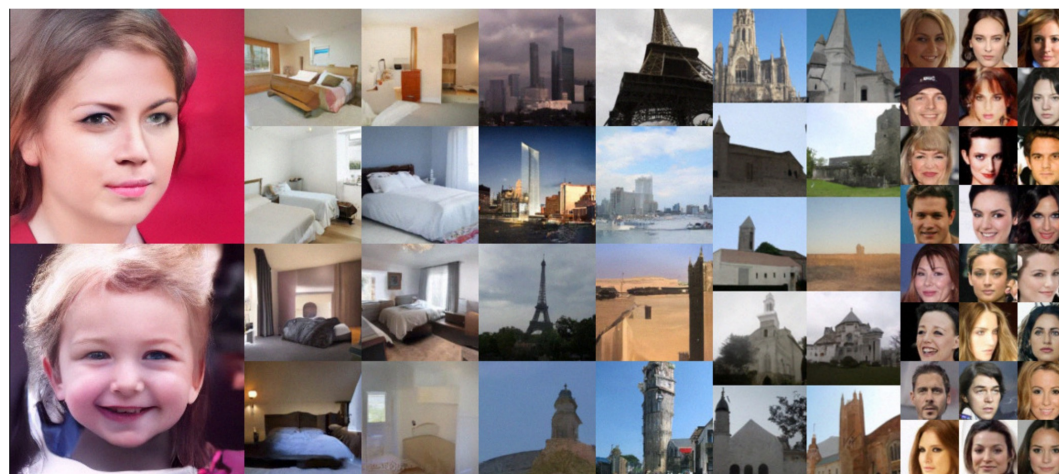


Figure 1: Generated samples on datasets of decreasing resolutions. From left to right: FFHQ 256×256 , LSUN bedroom 128×128 , LSUN tower 128×128 , LSUN church_outdoor 96×96 , and CelebA 64×64 .

EBMs (BP through CNNs) \rightarrow Score-based models (U-Nets)

[3] Song & Ermon. Generative Modeling by Estimating Gradients of the Data Distribution. NeurIPS 2019

[4] Song & Ermon. Improved Techniques for Training Score-Based Generative Models. NeurIPS 2020

Diffusion Models in 2021 (Stochastic Differential Equations)

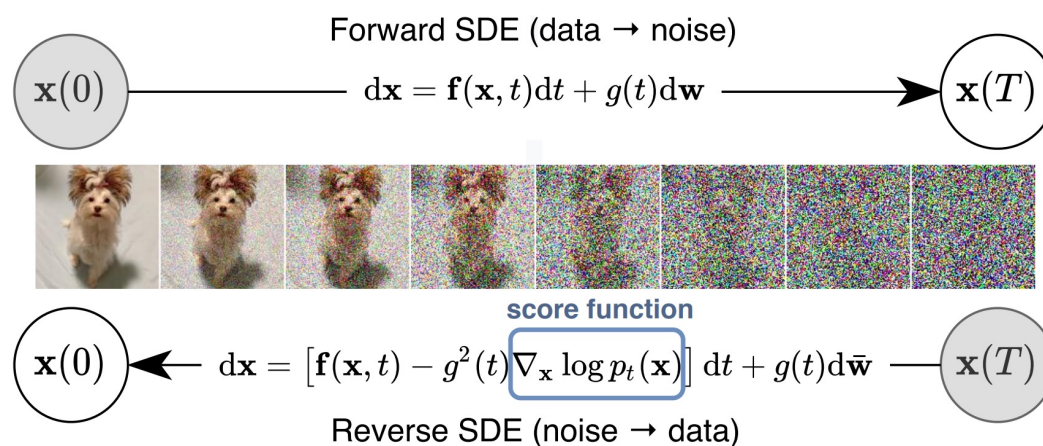


Figure 1: **Solving a reverse-time SDE yields a score-based generative model.** Transforming data to a simple noise distribution can be accomplished with a continuous-time SDE. This SDE can be reversed if we know the score of the distribution at each intermediate time step, $\nabla_x \log p_t(\mathbf{x})$.

- Drift coefficient f
- Diffusion coefficient g

Diffusion Processes



Forward process (transition distribution):

$$x_0 \sim q_0(x_0), \quad q_{0t}(x_t|x_0) = \mathcal{N}(x_t|\alpha_t x_0, \sigma_t^2 \mathbf{I})$$

Forward process (SDE):

$$dx_t = f(t)x_t dt + g(t)d\omega_t$$

$$\text{where } f(t) = \frac{d \log \alpha_t}{dt} \text{ and } g(t)^2 = \frac{d\sigma_t^2}{dt} - 2 \frac{d \log \alpha_t}{dt} \sigma_t^2$$

Diffusion Processes



Reverse process (SDE):

$$dx_t = \left[f(t)x_t - g(t)^2 \nabla_{x_t} \log q_t(x_t) \right] dt + g(t) d\bar{\omega}_t$$

Reverse process (ODE):

$$\frac{dx_t}{dt} = f(t)x_t - \frac{1}{2}g(t)^2 \nabla_{x_t} \log q_t(x_t)$$

Training DPMs by Score Matching



$$\mathcal{J}_{\text{SM}}^t(\theta) \triangleq \frac{1}{2} \mathbb{E}_{q_t(x_t)} \left[\left\| \mathbf{s}_{\theta}^t(x_t) - \nabla_{x_t} \log q_t(x_t) \right\|_2^2 \right]$$

Unknown

$$\mathcal{J}_{\text{SM}}(\theta; \lambda(t)) \triangleq \int_0^T \lambda(t) \mathcal{J}_{\text{SM}}^t(\theta) dt$$

Training DPMs by Denoising Score Matching



$$\mathcal{J}_{\text{DSM}}^t(\theta) \triangleq \frac{1}{2} \mathbb{E}_{q_0(x_0), q(\epsilon)} \left[\left\| \mathbf{s}_{\theta}^t(x_t) + \frac{\epsilon}{\sigma_t} \right\|_2^2 \right]$$

where $x_t = \alpha_t x_0 + \sigma_t \epsilon$ and $q(\epsilon) = \mathcal{N}(\epsilon | \mathbf{0}, \mathbf{I})$

The Stochastic Process of Data Score is a Martingale



Theorem 1. (Proof in Appendix A.1) Let $q_t(x_t)$ be constructed from the forward process in Eq. (2). Then under some regularity conditions, we have $\forall 0 \leq s < t \leq T$,

$$\alpha_t \nabla_{x_t} \log q_t(x_t) = \mathbb{E}_{q_{st}(x_s|x_t)} [\alpha_s \nabla_{x_s} \log q_s(x_s)], \quad (6)$$

where $q_{st}(x_s|x_t) = \frac{q_{st}(x_t|x_s)q_s(x_s)}{q_t(x_t)}$ is the transition probability from x_t to x_s .

Leads to concentration bounds and naturally $\mathbb{E}_{q_t(x_t)} [\nabla_{x_t} \log q_t(x_t)] = 0$

Calibrating DPMs



Although $\mathbb{E}_{q_t(x_t)} [\nabla_{x_t} \log q_t(x_t)] = 0$

Typically there is $\mathbb{E}_{q_t(x_t)} [\mathbf{s}_\theta^t(x_t)] \neq 0$

So we calibrate DPMs into $\mathbf{s}_\theta^t(x_t) - \eta_t$

Calibrating DPMs



Given any pretrained DPM, we can calibrate it as:

$$\mathbf{s}_\theta^t(x_t) - \mathbb{E}_{q_t(x_t)} [\mathbf{s}_\theta^t(x_t)]$$



$$\mathcal{J}_{\text{SM}}^t(\theta, \eta_t^*) = \mathcal{J}_{\text{SM}}^t(\theta) - \frac{1}{2} \left\| \mathbb{E}_{q_t(x_t)} [\mathbf{s}_\theta^t(x_t)] \right\|_2^2$$

Likelihood of Calibrating DPMs: SDE Solver



$$dx_t = [f(t)x_t - g(t)^2 \mathbf{s}_\theta^t(x_t)] dt + g(t) d\bar{\omega}_t$$



Marginal distribution

$$p_t^{\text{SDE}}(x_t; \theta)$$

$$\mathcal{D}_{\text{KL}}(q_0 \| p_0^{\text{SDE}}(\theta)) \leq \mathcal{J}_{\text{SM}}(\theta; g(t)^2) + \mathcal{D}_{\text{KL}}(q_T \| p_T)$$

Likelihood of Calibrating DPMs: SDE Solver



$$dx_t = \left[f(t)x_t - g(t)^2 (\mathbf{s}_\theta^t(x_t) - \underline{\eta}_t) \right] dt + g(t) d\bar{\omega}_t$$



Marginal distribution

$$p_0^{\text{SDE}}(x_t; \theta, \underline{\eta}_t)$$

$$\mathcal{D}_{\text{KL}}(q_0 \| p_0^{\text{SDE}}(\theta, \underline{\eta}_t)) \leq \mathcal{J}_{\text{SM}}(\theta, \underline{\eta}_t; g(t)^2) + \mathcal{D}_{\text{KL}}(q_T \| p_T)$$



Likelihood of Calibrating DPMs: *SDE Solver*

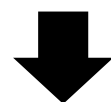
$$\mathcal{J}_{\text{SM}}(\theta, \eta_t^*; g(t)^2) = \mathcal{J}_{\text{SM}}(\theta; g(t)^2) - \frac{1}{2} \int_0^T g(t)^2 \left\| \mathbb{E}_{q_t(x_t)} [\mathbf{s}_\theta^t(x_t)] \right\|_2^2 dt$$

Upper bound reduced by calibration

Likelihood of Calibrating DPMs: ODE Solver



$$\frac{dx_t}{dt} = f(t)x_t - \frac{1}{2}g(t)^2 \mathbf{s}_\theta^t(x_t)$$



Marginal distribution

$$p_t^{\text{ODE}}(x_t; \theta)$$

$$\mathcal{D}_{\text{KL}}(q_0 \| p_0^{\text{ODE}}(\theta)) \approx \mathcal{J}_{\text{SM}}(\theta; g(t)^2) + \mathcal{D}_{\text{KL}}(q_T \| p_T)$$

Likelihood of Calibrating DPMs: ODE Solver



$$\frac{dx_t}{dt} = f(t)x_t - \frac{1}{2}g(t)^2(\mathbf{s}_\theta^t(x_t) - \underline{\eta_t})$$

↓ Marginal distribution

$$p_t^{\text{ODE}}(x_t; \theta, \underline{\eta_t})$$

$$\mathcal{D}_{\text{KL}}(q_0 \| p_0^{\text{ODE}}(\theta, \underline{\eta_t})) \approx \mathcal{J}_{\text{SM}}(\theta, \underline{\eta_t}; g(t)^2) + \mathcal{D}_{\text{KL}}(q_T \| p_T)$$

Empirical Results

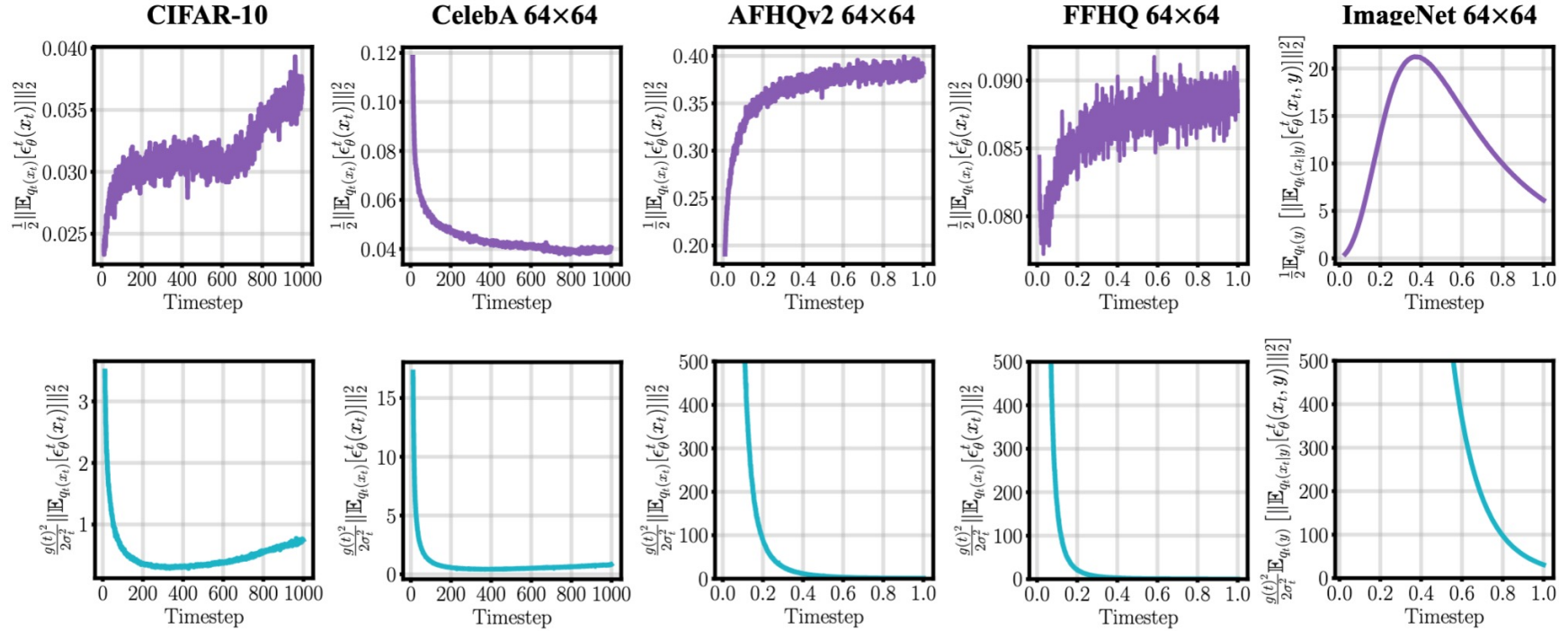


Figure 1: Time-dependent values of $\frac{1}{2} \|\mathbb{E}_{q_t(x_t)}[\epsilon_\theta^t(x_t)]\|_2^2$ (the first row) and $\frac{g(t)^2}{2\sigma_t^2} \|\mathbb{E}_{q_t(x_t)}[\epsilon_\theta^t(x_t)]\|_2^2$ (the second row) calculated on different datasets. The models on CIFAR-10 and CelebA is trained on discrete timesteps ($t = 0, 1, \dots, 1000$), while those on AFHQv2, FFHQ, and ImageNet are trained on continuous timesteps ($t \in [0, 1]$). We convert data prediction $\hat{x}_\theta^t(x_t)$ into noise prediction $\epsilon_\theta^t(x_t)$ based on $\epsilon_\theta^t(x_t) = (x_t - \alpha_t \hat{x}_\theta^t(x_t))/\sigma_t$. The y-axis is clamped into $[0, 500]$.

Empirical Results

Table 1: Comparison on sample quality measured by FID \downarrow with varying NFE on CIFAR-10. Experiments are conducted using a linear noise schedule on the discrete-time model from [15]. We consider three variants of DPM-Solver with different orders. The results with \dagger mean the actual NFE is $\text{order} \times \lfloor \frac{\text{NFE}}{\text{order}} \rfloor$ which is smaller than the given NFE, following the setting in [26].

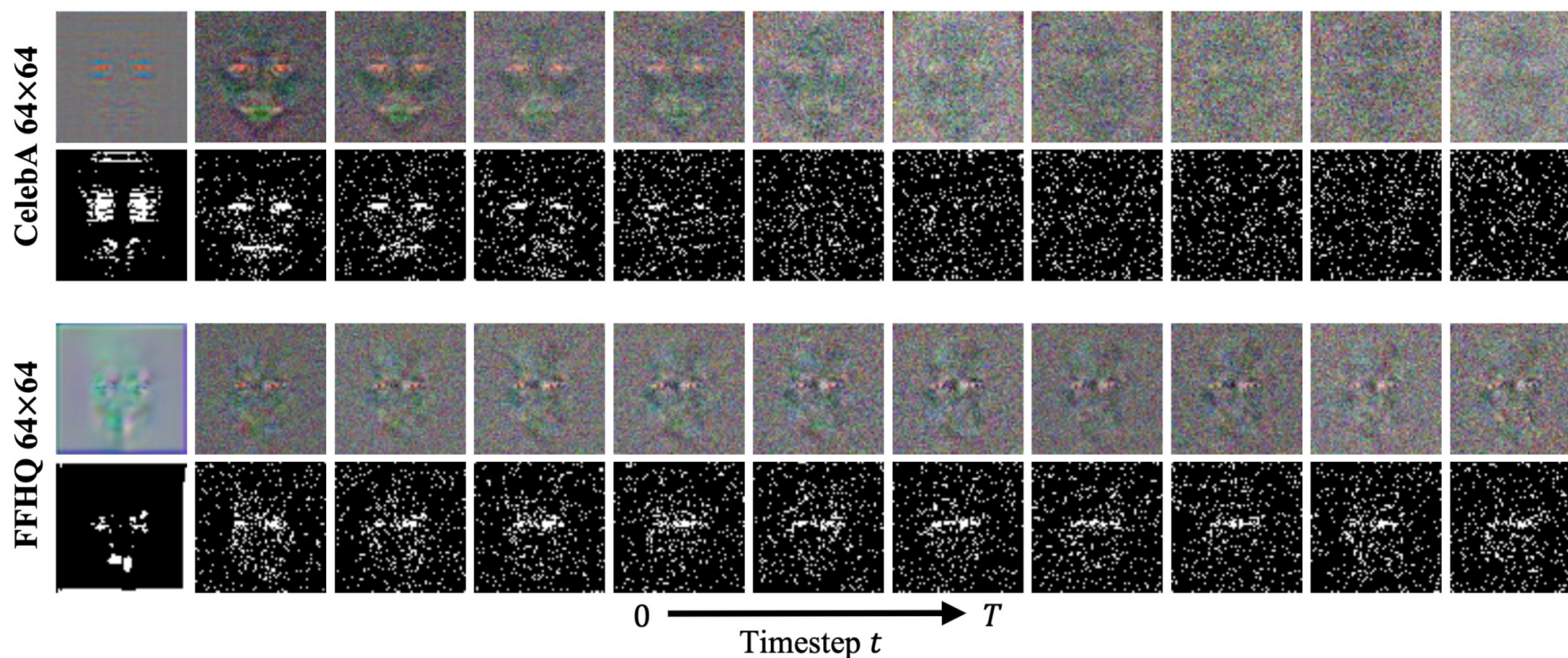
Noise prediction	DPM-Solver	Number of evaluations (NFE)						
		10	15	20	25	30	35	40
$\epsilon_{\theta}^t(x_t)$	1-order	20.49	12.47	9.72	7.89	6.84	6.22	5.75
	2-order	7.35	\dagger 4.52	4.14	\dagger 3.92	3.74	\dagger 3.71	3.68
	3-order	\dagger 23.96	4.61	\dagger 3.89	\dagger 3.73	3.65	\dagger 3.65	\dagger 3.60
$\epsilon_{\theta}^t(x_t) - \mathbb{E}_{q_t(x_t)}[\epsilon_{\theta}^t(x_t)]$	1-order	19.31	11.77	8.86	7.35	6.28	5.76	5.36
	2-order	6.76	\dagger 4.36	4.03	\dagger 3.66	3.54	\dagger 3.44	3.48
	3-order	\dagger 53.50	4.22	\dagger 3.32	\dagger 3.33	3.35	\dagger 3.32	\dagger 3.31

Table 2: Comparison on sample quality measured by FID \downarrow with varying NFE on CelebA 64×64 . Experiments are conducted using a linear noise schedule on the discrete-time model from [35]. The settings of DPM-Solver are the same as on CIFAR-10.

Noise prediction	DPM-Solver	Number of evaluations (NFE)						
		10	15	20	25	30	35	40
$\epsilon_{\theta}^t(x_t)$	1-order	16.74	11.85	7.93	6.67	5.90	5.38	5.01
	2-order	4.32	\dagger 3.98	2.94	\dagger 2.88	2.88	\dagger 2.88	2.84
	3-order	\dagger 11.92	3.91	\dagger 2.84	\dagger 2.76	2.82	\dagger 2.81	\dagger 2.85
$\epsilon_{\theta}^t(x_t) - \mathbb{E}_{q_t(x_t)}[\epsilon_{\theta}^t(x_t)]$	1-order	16.13	11.29	7.09	6.06	5.28	4.87	4.39
	2-order	4.42	\dagger 3.94	2.61	\dagger 2.66	2.54	\dagger 2.52	2.49
	3-order	\dagger 35.47	3.62	\dagger 2.33	\dagger 2.43	2.40	\dagger 2.43	\dagger 2.49



Empirical Results



Thanks

